

# Data Analysis and Robust Modelling of the Impact of Renewable Generation on Long Term Security of Supply and Demand

Matthias C. M. Troffaes

Dept. of Mathematical Sciences  
Durham University, UK

Email: matthias.troffaes@gmail.com

Edward Williams

Dept. of Mathematical Sciences  
Durham University, UK

Email: e.h.williams@durham.ac.uk

Chris Dent

School of Engineering & Computing Sciences  
Durham University, UK

Email: chris.dent@durham.ac.uk

**Abstract**—This paper studies rigorous statistical techniques for modelling long term reliability of demand and supply of electrical power given uncertain variability in the generation and availability of wind power and conventional generation. In doing so, we take care to validate statistical assumptions, using historical observations, as well as our intuition about the actual underlying real-world statistical process. Where assumptions could not be easily validated, we say so explicitly. In particular, we aim to improve existing statistical models through sensitivity analysis of ill-known parameters: we propose models for wind power and conventional generation, estimate their parameters from historical wind power data and conventional availability data, and finally combine them with historical demand data to build a full robust joint time-dependent model of energy not served. Bounds on some useful indices from this model are then calculated, such as expected energy not served, and expected number of continuous outage periods—the latter cannot be estimated from a purely time collapsed model because time collapsed models necessarily do not model correlations across time. We compare our careful model with a naive model that ignores deviations from normality, and find that this results in substantial differences: in this specific study, the naive model overestimates the risk roughly by a factor 2. This justifies the care and caution by which model assumptions must be verified, and the effort that must be taken to adapt the model accordingly.

## I. INTRODUCTION

The inclusion of variable generation within power system adequacy risk calculations is currently a key topic in power system planning methodology. A vital component of this is an appropriate statistical wind resource model. Some outputs of interest, such as for example the expected energy not served (which is the expected total energy shortage in a future time window), can be calculated using a time collapsed model in which time correlations in variable generation are not modelled explicitly [1]. However, other outputs of interest, such as for instance the expected number of periods of shortfall, require a full time series model of the variable generation. Examples of statistical approaches to this may be found in [2], [3], [4].

This paper makes two contributions in the use of time series wind models within power system adequacy calculations. First, although standard ARMA processes have Gaussian marginal distributions, it is not standard practice in the power system

literature to transform wind speed data or wind power data so that it has a Gaussian marginal before estimating parameters of an ARMA wind speed or wind power model—although there are some exceptions [5], [6]. We will demonstrate that ensuring the conditions for an ARMA process are correctly satisfied can make a substantial difference to model outputs. Secondly, as we shall see, the parameters of the ARMA process can vary substantially when fitted to data from different years. We demonstrate that these differences can lead to quite different results, and propose a method for sensitivity analysis, based on imprecise probability [7], [8].

The analysis is based on the “Adjusted Gone Green” scenario supplied by National Grid, in which the generating unit capacities are slightly adjusted from the original “Gone Green” scenario. The results are thus generally representative of Great Britain calculations and are entirely sufficient for demonstrating methodology. The small data adjustments are necessary in order to make clear that model outputs do not precisely reflect any future scenario for the Great Britain system developed by National Grid.

Sections II and III discuss, respectively, the wind power model and the conventional generation model. Various risk indices are derived from these models in section IV. We reflect on the results and future work in section V.

## II. MODELLING WIND POWER

The available wind power data covers seven twenty week winters. The aim of our model is to characterize the statistical properties of the wind power time series in an arbitrary year.

A simple and often effective approach to time series modelling is to use an ARMA process [9]. ARMA processes have a normal marginal distribution, but obviously our wind power will not be normally distributed.

The top left plot in fig. 1 confirms that the data is not normally distributed; there are fewer values in the tail of the distribution than there would be in a normal distribution. Also, if other years are to be simulated from the model, there are certain constraints with which simulated data has to abide; no output can be below 0 or above 10120 which is the maximum wind power output of the scenario on which the data is based.

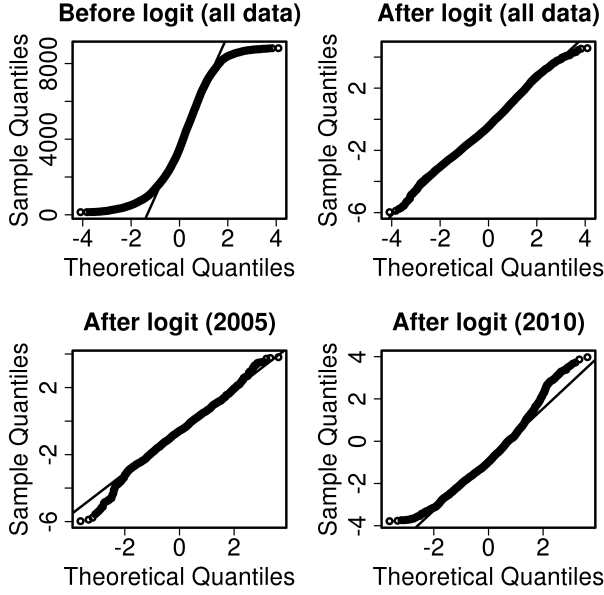


Fig. 1. Normal quantile-quantile plots for the wind data.

In reality we can set this interval to be smaller than  $[0, 10120]$  as it is highly improbable to observe an output at extremes.

Various approaches to dealing with non-normality in wind power data are discussed in the literature; see for instance [5], [6]. A simple solution is to take a logit transform:

$$\text{logit}(x) = \log\left(\frac{t(x)}{1-t(x)}\right) \text{ where } t(x) = \frac{x-\alpha}{\beta-\alpha}. \quad (1)$$

Here,  $[\alpha, \beta]$  represents the range of possible values of  $y$ .

Choosing  $\alpha = 0$  and  $\beta = 10120$  according to the physical constraints of the system does not entirely fix departures from normality. By visual inspection of quantile-quantile plots,  $\alpha = 120$  and  $\beta = 8900$  were selected. The quantile-quantile plot for this choice is depicted in the top right of fig. 1; clearly there is significant improvement. It is important to realize that our model thereby excludes extreme wind power events outside  $[120, 8900]$ . Note that the actual bounds of full 7 years of wind power data are  $[142.306, 8810.187]$ , so observations outside the range  $[120, 8900]$  were never actually observed in the data.

It is instructive to inspect also the quantile-quantile plots for individual years. The most problematic years are 2005 and 2010, depicted in the bottom of fig. 1, with the actual quantiles deviating from the theoretical quantiles in the upper tail for 2010 and the lower tail for 2005. These observations may suggest that the bounds are not entirely static, and perhaps some random effect on  $\alpha$  and  $\beta$  could be included. For simplicity, however, we will stick with constant bounds.

We propose the following model:

$$\text{logit}(X(y, t)) = Z_1(y) + Z_2(y, t) \quad (2)$$

where  $y$  is the year,  $t$  is the time within the year,  $Z_1(y)$  captures a yearly effect, and  $Z_2(y, t)$  is an ARMA process with zero mean. Because

$$E(\text{logit}(X(y, t))|y) = E(Z_1(y)|y) + 0 = Z_1(y) \quad (3)$$

TABLE I  
ESTIMATED REALIZATIONS OF  $Z_1(y)$ , WITH CONFIDENCE LIMITS AT THE 95% CONFIDENCE LEVEL.

$y$	2005	2006	2007	2008	2009	2010	2011
$\hat{z}_1(y)$	-0.62	0.22	0.03	-0.56	-0.66	-0.84	-0.17
error	$\pm 0.33$	$\pm 0.56$	$\pm 0.42$	$\pm 0.38$	$\pm 0.45$	$\pm 0.32$	$\pm 0.61$

we can estimate  $Z_1(y)$  from the yearly sample mean of  $\text{logit}(X(y, t))$ . As we only have 7 years of observations, we can only estimate 7 realizations of  $Z_1(y)$ . The exact values thus estimated are listed in table I, along with confidence limits, where the autocorrelation was taken into account (for example, see [10, Sec. 3.1]). It can be seen that there is a significant difference between each yearly effect. However, most of the confidence intervals for  $\hat{z}_1(y)$  overlap with each other, making it difficult to draw any strong conclusions.

Because the entire process is reasonably normal, we judge it not unreasonable to assume that the  $Z_1(y)$  are iid samples from a normal distribution:

$$Z_1(y) \sim N(\mu, \sigma^2) \quad (4)$$

where we can estimate  $\mu$  and  $\sigma^2$  by

$$\hat{\mu} = \frac{1}{7} \sum_{y=1}^7 \hat{z}_1(y) = -0.371 \quad (5)$$

$$\hat{\sigma}^2 = \frac{1}{6} \sum_{y=1}^7 (\hat{z}_1(y) - \hat{\mu})^2 = 0.397^2 \quad (6)$$

The error on  $\hat{\mu}$ , at the 95% level, is:

$$\begin{aligned} 1.96 \times \text{s.e.}(\hat{\mu}) &= 1.96 \left( \underbrace{\frac{\hat{\sigma}^2}{7}}_{0.022} + \underbrace{\frac{1}{7^2} \sum_{y=1}^7 \text{var}(\hat{z}_1(y))}_{0.008} \right)^{0.5} \\ &= 0.34 \end{aligned} \quad (8)$$

The errors in  $\hat{z}_1(y)$  are the least contribution to the total error. Ignoring the errors in  $\hat{z}_1(y)$ , a naive but simple 95% confidence interval for  $\hat{\sigma}^2$  is:

$$[6\hat{\sigma}^2/\chi_{0.975}^2, 6\hat{\sigma}^2/\chi_{0.025}^2] = [0.256^2, 0.874^2] \quad (9)$$

where  $\chi_{\alpha}^2$  are the  $\alpha$  quantiles of the chi-square distribution with six degrees of freedom.

When simulating a random year, we may simply draw the yearly effect from  $N(\hat{\mu}, \hat{\sigma}^2)$ . Alternatively, a more conservative analysis could simply take the lowest value for  $Z_1(y)$  within the confidence intervals of table I, that is,  $z_1 = -0.84 - 0.32 = -1.16$ . This could be appropriate if one would like to drop the assumption that the  $Z_1(y)$  are iid realisations from a normal distribution.

As said before, we will assume that  $Z_2$  is an ARMA process with zero mean. More specifically, we assume that

$$Z_2(y, \bullet)|y \sim \text{AR}(\alpha(y), \sigma(y)) \quad (10)$$

where  $\alpha(y)$  are the coefficients of the AR process and  $\sigma(y)$  is the standard deviation of the residual noise—these parameters may vary across years. An estimate for  $Z_2(y, t)$  is

$$\hat{z}_2(y, t) = \text{logit}(x(y, t)) - \hat{z}_1(y). \quad (11)$$

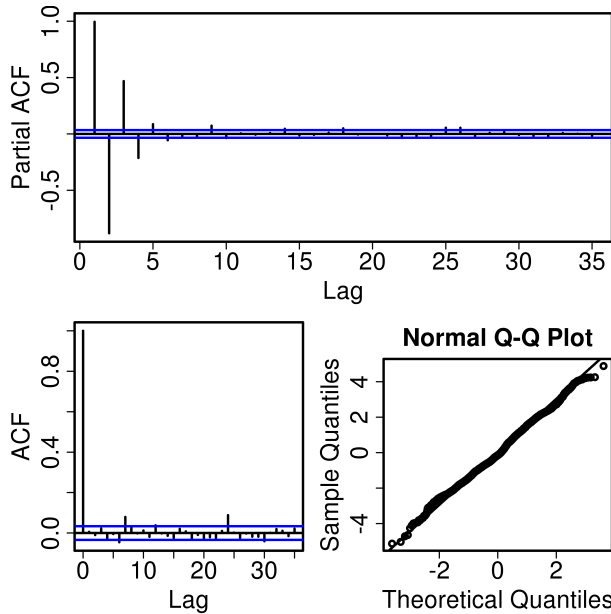


Fig. 2. Top: partial autocorrelation diagram for  $Z_2(y, t)$ , for  $y = 2008$ . Bottom: full autocorrelation diagram and normal quantile-quantile plot for the residuals of the fitted  $AR(5)$  model for  $Z_2(y, t)$ , for  $y = 2008$ .

To judge the stationarity of  $Z_2(y, t)$ , and to pick an appropriate order for the AR process, we investigate the partial autocorrelation diagram of the data, depicted for  $y = 2008$  in fig. 2. Other years follow a very similar pattern. Larger autocorrelations are limited, thus a stationarity assumption seems reasonable. The partial autocorrelation become close to zero at a lag of 4 or 5, which would suggest an AR model of order 4 or 5. Therefore, we used an AR model of order 5.

To check the model fit, the full autocorrelation diagram and normal quantile-quantile plot of the residual of the fitted model are also plotted in fig. 2, again for  $y = 2008$ . There are no large peaks left, so the model is adequate as the residuals do not have any significant correlation. The largest peak occurs at lag 24, which corresponds to a day. This suggests that there might be a daily cycle that has not been taken into account. Because the peak is quite small, no further attempt to remove it was made in this study.

To finalise the model, we need to fit the AR model coefficients. Naively, we could join all of the years together after we have transformed each year and set its mean to 0, hoping that all years are similar enough, and hoping that the discontinuity across years has little impact on our estimates. More cautiously, we could fit an AR model to each of the years separately. As can be seen from table II, the coefficients vary significantly across years. The errors on these estimates are consistently approximately equal to

$$1.96 \times \text{s.e.}(\hat{\alpha}(y)) \approx (0.03, 0.09, 0.11, 0.09, 0.03) \quad (12)$$

across all years  $y$ . Clearly, the variation of the estimated  $\alpha(y)$  across years is far larger than the errors on the estimates. To proceed cautiously, we merely assume that at least one of the

TABLE II  
FITTED AR COEFFICIENTS FOR EACH YEAR.

$y$	$\hat{\alpha}_1(y)$	$\hat{\alpha}_2(y)$	$\hat{\alpha}_3(y)$	$\hat{\alpha}_4(y)$	$\hat{\alpha}_5(y)$	$\hat{\sigma}(y)$
2005	2.54	-2.54	1.48	-0.64	0.16	0.04
2006	2.56	-2.65	1.63	-0.7	0.16	0.06
2007	2.49	-2.45	1.38	-0.55	0.12	0.06
2008	2.41	-2.25	1.17	-0.44	0.09	0.06
2009	2.56	-2.58	1.46	-0.58	0.14	0.04
2010	2.53	-2.51	1.39	-0.54	0.12	0.04
2011	2.22	-1.73	0.68	-0.27	0.09	0.06

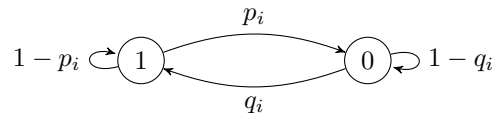


Fig. 3. Two-state Markov chain for conventional capacity.

years is representative for a future year, however we do not know which of the years. In particular, we do not assume that the  $Z_2(y, \cdot)$  processes are fully exchangeable across years  $y$ . For inference, we will simply do a sensitivity analysis on the coefficients from each observed year, and bound the resulting probabilities and expectations [7], [8].

### III. MODELLING CONVENTIONAL GENERATION

For each unit of conventional generation, we have its capacity  $c_i$ , and the fraction of time  $a_i$  it is available. To model the total conventional capacity in time, we assume that each unit  $W_i$  follows a 2 state discrete time Markov chain, where at each time point  $t$  the unit is either working ( $W_i(t) = 1$ ) or not working ( $W_i(t) = 0$ ); see fig. 3. Time steps by the hour. An hourly resolution provides sufficient detail for our purpose. Moreover, our wind power data is also by the hour.

Each unit  $i$  then has two parameters:  $p_i$  and  $q_i$ . The mean time to repair is  $1/q_i$ . Due to lack of data, we assume identical repair rates across all conventional generators. Because 50 hours mean time to repair is reasonably representative of typical generation units [11, Table 1], we simply set  $q_i = 1/50$  for all units  $i$ —obviously this aspect of the model could be improved in future work. The theoretical long term availability is simply the limiting probability of  $W_i = 1$ , which is equal to  $\frac{q_i}{p_i + q_i}$ . Consequently, the equality  $a_i = \frac{q_i}{p_i + q_i}$  determines  $p_i$ , as  $a_i$  and  $q_i$  are known.

The total available capacity for conventional is then simply

$$X(t) = \sum_{i=1}^k c_i W_i(t). \quad (13)$$

For simulating  $X(t)$ , we bluntly simulate  $W_i(t)$  for each conventional unit and then join these together using eq. (13).

### IV. ENERGY NOT SERVED AND NUMBER OF SHORTFALLS

The *energy not served* is defined as

$$E := \sum_{t=1}^{3360} \max\{0, D(t) - C(t) - W(t)\} \quad (14)$$

where  $D(t)$ ,  $C(t)$ , and  $W(t)$  are, respectively, the demand, conventional generation, and wind generation at time  $t$ . The

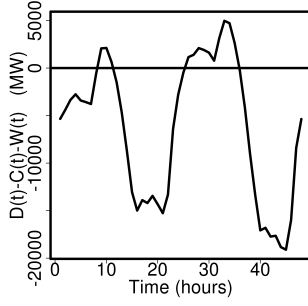


Fig. 4. The energy not served  $E$  is the area under the curve which lies above the horizontal axis. The number of shortfalls  $N$  is the number of such areas.

sum runs over 3360 hours, which is equal to the length of each winter period in the data. The *number of shortfalls*,  $N$ , is defined as the number of times that the sequence  $(D(t) - C(t) - W(t))_{t=1 \dots 3360}$  changes sign from negative to positive.

Figure 4 demonstrates both concepts, for a shorter period of 48 hours. Note that, to serve clarity, the plot is schematic and does not represent results from our modelling. Anyway, in this case,  $E = 29251$  and  $N = 2$ .

Our main aim is to estimate the distributions of  $E$  and  $N$ , and thereby also various indices of system reliability, such as the expectation of  $E$ , that is, the expected energy not served, often abbreviated as EENS. Note that the distribution of  $N$  cannot be estimated from a purely time collapsed model because time collapsed models necessarily do not model correlations across time. To do this estimation, we combine the simulated wind and simulated conventional capacity along with historic demand data. Demand simulation is not attempted, as this requires complicated modelling of periodic effects at many different time scales, which we have not yet attempted to model. Instead, we simply use an actual demand trace from 2010, which has the highest peak demand across all years between 2005 and 2011. Figure 5 shows a short simulation trace of wind and conventional, along with demand.

In a normal situation, we would perform  $n$  model runs, where each model run simulates a full winter (3360 hours) of wind generation and conventional capacity, thereby producing a single realisation of  $E$  and  $N$ . However, as discussed in section II, we fitted 7 distinct AR models for  $Z_2(y, \cdot)$ , namely one for each winter in the data. To allow us to perform a sensitivity analysis against the AR model parameters, every model run simulated wind for each of these AR models, thereby producing 7 distinct realisations of  $E$  and  $N$ . To ensure consistent sampling errors, the random seed of each of the 7 wind simulations was forced to the same value—of course this seed varied randomly between simulation runs.

We then use these samples of  $E$  and  $N$  to produce lower and upper expectations, or more sophisticatedly, lower and upper histograms. Our actual simulations used  $n = 100\,000$ , and the total time to complete these runs was 114.21 hours, with most time spent on simulation of conventional capacity.

The top left of fig. 6 shows 7 overlaid histograms—

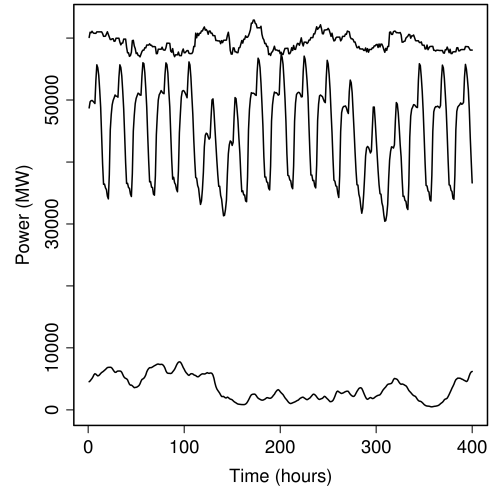


Fig. 5. A trace of wind generation (bottom curve), historic demand (middle curve), and conventional capacity (top curve).

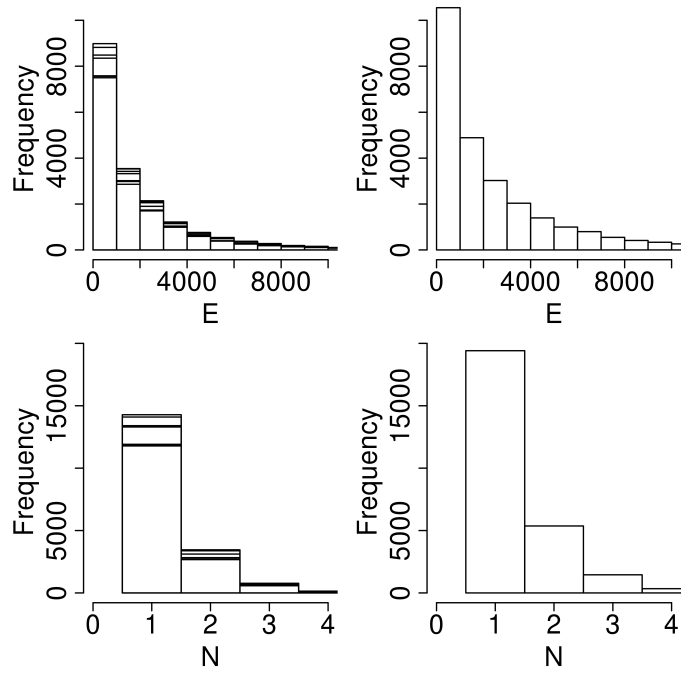


Fig. 6. Estimated distributions of  $E \mid E > 0$  and  $N \mid N \geq 1$ . The histograms on the left show the model with logit and year effect for wind. The histograms on the right show a naively fitted model without logit and without year effect.

one histogram for every choice of AR parameters—for  $E$ , conditional on  $E > 0$ ; the large peak at  $E = 0$  (see eq. (19) further) has been omitted to make for a clearer picture. The largest simulated value for  $E$  was 43653.5, however  $E \geq 10\,000$  is very rare (only a fraction 0.00412 of all model runs), so this tail has been omitted from the histogram. The top right of fig. 6 shows the results of a more naive model for wind, without logit and without random year effect (or more precisely, a constant year effect equal to the mean of the data). In this case, this has led to a clear overestimation of the risk.

Similar histograms for  $N$  are depicted at the bottom of fig. 6. Again,  $N = 0$  is omitted for clarity. The largest simulated value for  $N$  was 8, however  $N \geq 4$  is very rare (only a fraction 0.00148 of all model runs), so this tail has been omitted. If there is a shortfall in a year, then it is likely to have only happened once or twice. The chance of having a shortfall more than twice in a single year, given that shortfall occurs, is at most 0.054. Obviously, here too, a naive model overestimates the risk quite substantially.

Lower and upper expected energy not served, that is, bounds on the EENS, can be estimated from the simulation as follows:

$$\underline{P}(E) = \min_{y=1}^7 \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{3360} \max\{0, d(t) - c_i(t) - w_{yi}(t)\} \quad (15)$$

$$= 299.63 \pm 9.24 \quad (16)$$

$$\overline{P}(E) = \max_{y=1}^7 \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{3360} \max\{0, d(t) - c_i(t) - w_{yi}(t)\} \quad (17)$$

$$= 389.91 \pm 9.24 \quad (18)$$

where the error is the worst case 95% confidence interval around the sample means,  $n = 100\,000$  is the number of model runs,  $d(t)$  is historic demand at time  $t$  (for 2010),  $c_i(t)$  is conventional capacity at time  $t$  in model run  $i$ , and  $w_{yi}(t)$  is the wind power at time  $t$  in model run  $i$  with AR coefficients  $\hat{\alpha}(y)$ . Bounds on the probability of  $E = 0$  follow similarly:

$$\underline{P}(E = 0) = 0.813 \pm 0.002 \quad \overline{P}(E = 0) = 0.848 \pm 0.002 \quad (19)$$

For comparison, the naive model has:

$$P(E) = 808.416 \pm 16.501 \quad P(E = 0) = 0.733 \pm 0.003 \quad (20)$$

which confirms our earlier observation about risk overestimation. In particular, the naive model overestimates the expected energy not served by a factor of more than 2.

The difference between AR coefficients  $\hat{\alpha}(y)$  clearly has a significant impact, as do statistical assumptions such as normality, justifying the caution by which we fitted our models.

## V. CONCLUSION

We modelled wind and conventional power to characterize the distributions of typical quantities of interest, such as energy not served, and number of shortfalls. Model parameters varied substantially when fitted to data from different years. In our analysis, these differences lead to quite different results, prompting a sensitivity analysis. We showed that a naive model—ignoring non-normality and differences across years—leads to substantial overestimation of the risk. Even though the model in this study was quite limited, it does highlight the importance of careful statistical modelling.

Various aspects of the model remain to be improved. The distributions of energy not served and number of shortfalls would be as in this paper only if a future winter had the same demand pattern, if conventional capacity actually followed a two state Markov chain, and if the wind was similar to wind seen in at least one of the winters in our data. Therefore, we have not made a strong ‘real world’ statement.

First, we used a simple historic trace for demand. Due to complex periodic effects at different time scales (daily, weekends, season) and dependence on climate and economy, predictive modelling of demand is a non-obvious task [12].

Next, our model for conventional generation has quite a few limitations. In particular, the Markovian assumption may be violated, repair rates will not be equal across all conventional generators, and repair and failure rates will also not be constant throughout the day: generators fail more regularly during startup and during ramping.

Finally, the logit transform, which was used to transform our wind into a normally distributed process, implies that we cannot simulate wind power output larger than the bounds on the transformation, which were set quite closely to the maximum and minimum value observed from the data. Perhaps these bounds should not be taken to be constant.

## ACKNOWLEDGMENT

We are grateful to National Grid for providing the data that was used in this study. The research reported in this paper was funded by BP through the Durham Energy Institute, and by EPSRC (grant no EP/K002252/1).

## REFERENCES

- [1] C. J. Dent and S. Zachary, “Further results on the probability theory of capacity value of additional generation,” in *International Conference on Probabilistic Methods Applied to Power Systems (PMAPS)*, Durham, England, 2014.
- [2] R. Billinton, Y. Gao, and R. Karki, “Composite system adequacy assessment incorporating large-scale wind energy conversion systems considering wind speed correlation,” *IEEE Transactions on Power Systems*, vol. 24, no. 3, pp. 1375–1382, Aug. 2009.
- [3] W. Wangde and R. Billinton, “Considering load-carrying capability and wind speed correlation of WECS in generation adequacy assessment,” *IEEE Transactions on Energy Conversion*, vol. 21, no. 3, pp. 734–741, Sep. 2006.
- [4] R. Karki, P. Hu, and R. Billinton, “A simplified wind power generation model for reliability evaluation,” *IEEE Transactions on Energy Conversion*, vol. 21, no. 2, pp. 533–540, Jun. 2006.
- [5] P. Pinson, “Very-short-term probabilistic forecasting of wind power with generalized logitnormal distributions,” *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, vol. 61, no. 4, pp. 555–576, 2012.
- [6] J. Tastu, P. Pinson, P.-J. Trombe, and H. Madsen, “Probabilistic forecasts of wind power generation accounting for geographically dispersed information,” *IEEE Transactions on Smart Grid*, vol. 5, no. 1, pp. 480–489, Jan. 2014.
- [7] P. Walley, *Statistical Reasoning with Imprecise Probabilities*. London: Chapman and Hall, 1991.
- [8] M. C. M. Troffaes and G. de Cooman, *Lower Previsions*, ser. Wiley Series in Probability and Statistics. Wiley, 2014. [Online]. Available: <http://eu.wiley.com/WileyCDA/WileyTitle/productCd-0470723777.html>
- [9] M. Falk, F. Marohn, R. Michel, D. Hofmann, M. Macke, C. Spachmann, and S. Englert, *A First Course on Time Series Analysis with SAS*. epubli GmbH, Aug. 2012. [Online]. Available: <http://statistik.mathematik.uni-wuerzburg.de/timeseries/>
- [10] C. J. Geyer, “Practical Markov chain Monte Carlo,” *Statistical Science*, vol. 7, no. 4, pp. 473–483, 1992. [Online]. Available: <http://www.jstor.org/stable/2246094>
- [11] A. B. Attya, Y. G. Hegazy, and M. A. Moustafa, “Random operation of conventional distributed generators based on generation techniques,” in *Canadian Conference on Electrical and Computer Engineering (CCECE 2008)*, May 2008, pp. 1203–1206.
- [12] C.-L. Hor, S. Watson, and S. Majithia, “Analyzing the impact of weather variables on monthly electricity demand,” *IEEE Transactions on Power Systems*, vol. 20, no. 4, pp. 2078–2085, Nov. 2005.